

PNAS Classification: Biochemistry

Title: Amino acid substitution matrices from protein blocks

Authors: Steven Henikoff<sup>\*,+</sup> and Jorja G. Henikoff

\*Howard Hughes Medical Institute  
Basic Sciences Division  
Fred Hutchinson Cancer Research Center  
Seattle, WA 98104  
Communicated by Walter Gilbert

Key Terms: amino acid sequence, alignment algorithms, database searching

Abbreviations: PAM, percent accepted mutation; BLOSUM, blocks substitution matrix

<sup>+</sup>To whom reprint requests should be addressed

**ABSTRACT** Methods for alignment of protein sequences typically measure similarity by using a substitution matrix with scores for all possible exchanges of one amino acid with another. The most widely-used matrices are based on the Dayhoff model of evolutionary rates. Using a different approach, we have derived substitution matrices from about 2000 blocks of aligned sequence segments characterizing more than 500 groups of related proteins. This led to marked improvements in alignments and in searches using queries from each of the groups.

Among the most useful computer-based tools in modern biology are those that involve sequence alignments of proteins, since these alignments often provide important insights into gene and protein function. There are several different types of alignments: "global" alignments of pairs of proteins related by common ancestry throughout their lengths, "local" alignments involving related segments of proteins, multiple alignments of members of protein families, and alignments made during database searches to detect homology. In each case, competing alignments are evaluated using a scoring scheme for estimating similarity. Although several different scoring schemes have been proposed (1-6), the mutation data matrices of Dayhoff (1, 7-9) are generally considered the standard and are often the default in alignment and searching programs. In the Dayhoff model, substitution rates are derived from alignments of protein sequences that are at least 85% identical. However, the most common task involving substitution matrices is the detection of much more distant relationships, which are only inferred from substitution rates in the Dayhoff model. Therefore, we wondered whether a better approach might be to use alignments in which these relationships are explicitly represented. An incentive for investigating this possibility is that implementation of an improved matrix in numerous important applications requires only trivial effort.

## METHODS

**Deriving a frequency table from a database of blocks.** Local alignments can be represented as ungapped "blocks" with each row a different protein segment and each column an aligned residue position. Previously, we described an automated system, PROTOMAT, for obtaining a set of blocks given a group of related proteins (10). This system was applied to a catalog of several hundred protein groups, yielding a database of over 2000 blocks. Consider a single block representing a conserved region of a protein family. For a new member of this family, we seek a set of scores for matches and mismatches that best favors a correct alignment with each of the other segments in the block relative to an incorrect alignment. For each column of the block, we first count the number of matches and mismatches of each type between the new sequence and every other sequence in the block. For example, if the residue of the new sequence that aligns with the first column of the first block is A and the column has 9 As and 1 S, then there are 9 AA matches and 1 AS mismatch. This procedure is repeated for all columns of all blocks with the summed results stored in a table. The new sequence is added to the group. For another new sequence the same procedure is followed, summing these numbers with those already in the table. Notice that successive addition of each sequence to the group leads to a table consisting of counts of all possible amino acid pairs in a column. For example, in the column consisting of 9 As and 1 S, there are  $8+7+\dots+1 = 36$  possible AA pairs, 9 AS or SA pairs and no SS pairs. Counts of all possible pairs in each column of each block in the database are summed. So, if a block has a width of  $w$  amino acids and a depth of  $s$  sequences, it contributes  $ws(s-1)/2$  amino acid pairs to the count ( $1 \times 10 \times 9 \div 2 = 45$  in the example above). The result of this counting is a frequency table listing the number of times each of the  $20+19+\dots+1 = 210$  different amino acid pairs occurs among the blocks. The table is used to calculate

a matrix representing the odds ratio between these observed frequencies and those expected by chance.

**Computing a log-odds matrix.** Let the total number of amino acid  $i,j$  pairs ( $1 \leq j \leq i \leq 20$ ) for each entry of the frequency table be  $f_{ij}$ . Then the observed probability of occurrence for each  $i,j$  pair is  
For the column of 9 As and 1 S in the example, where  $f_{AA} = 36$  and  $f_{AS} = 9$ ,  $q_{AA} =$

$$q_{ij} = f_{ij} / \sum \sum f_{ij}$$

$36 \div 45 = 0.8$  and  $q_{AS} = 9 \div 45 = 0.2$ . Next we estimate the expected probability of occurrence for each  $i,j$  pair. It is assumed that the observed pair frequencies are those of the population. For the example, 36 pairs have A in both positions of the pair and 9 pairs have A at only one of the two positions, so that the expected probability of A in a pair is  $[36 + (9 \div 2)] \div 45 = 0.9$  and of S is  $(9 \div 2) \div 45 = 0.1$ . In general, the probability of occurrence of the  $i$ th amino acid in an  $i,j$  pair is: The expected probability of occurrence  $e_{ij}$  for each  $i,j$  pair is then  $p_i p_j$  for  $i = j$  and

$$p_{ij} = q_{ii} + \sum q_{ij} / 2$$

$p_i p_j + p_j p_i = 2p_i p_j$  for  $i \neq j$ . In the example, the expected probability of AA is  $0.9 \times 0.9 = 0.81$ , of AS + SA is  $2 \times (0.9 \times 0.1) = 0.18$  and of SS is  $0.1 \times 0.1 = 0.01$ . An odds ratio matrix is calculated where each entry is  $q_{ij}/e_{ij}$ . A log-odds ratio is then calculated in bit units as  $s_{ij} = \log_2(q_{ij}/e_{ij})$ . If the observed frequencies are as expected,  $s_{ij} = 0$ , if less than expected,  $s_{ij} < 0$ , if more than expected,  $s_{ij} > 0$ . Log-odds ratios are multiplied by a scaling factor of 2 and then rounded to the nearest integer value to produce BLOSUM substitution matrices in half bit units, comparable to matrices generated by the PAM program (11). For each substitution matrix, we calculated the average mutual information (12) per amino acid pair  $H$  (also called relative entropy), and the expected score  $E$  in bit units as:

**Clustering segments within blocks.** To reduce multiple contributions to

$$H = \sum \sum q_{ij} \times s_{ij}, \quad E = \sum \sum p_i \times p_j \times s_{ij}.$$

amino acid pair frequencies from the most closely related members of a family, sequences are clustered within blocks and each cluster weighted as a single sequence in counting pairs (13). This is done by specifying a clustering percentage in which sequence segments that are identical for at least that percentage of amino acids are grouped together. For example, if the percentage is set at 80%, and sequence segment A is identical to sequence segment B at  $\geq 80\%$  of their aligned positions, then A and B are clustered and their contributions averaged in calculating pair frequencies. If C is identical to either A or B at  $\geq 80\%$  of aligned positions, it is also clustered with them and the contributions of A, B and C averaged, even though C might not be identical to both A and B at  $\geq 80\%$  of aligned positions. In the above example, if 8 of the 9 sequences with As in the 9A-1S column are clustered, then the contribution of this column to the frequency

table is equivalent to that of a 2A-1S column, which contributes 2 AS pairs. A consequence of clustering is that the contribution of closely related segments to the frequency table is reduced (or eliminated when an entire block is clustered, since this is equivalent to a single sequence in which no substitutions appear). For example, clustering at 62% reduces the number of blocks contributing to the table by 25%, with the remainder contributing 1.25 million pairs (including fractional pairs), whereas without clustering, more than 15 million pairs are counted (Fig. 1). In this way, varying the clustering percentage leads to a family of matrices. The matrix derived from a database of blocks in which sequence segments that are identical at  $\geq 80\%$  of aligned residues are clustered is referred to as BLOSUM 80, and so forth. The BLOSUM program implements matrix construction. Frequency tables, matrices and programs for UNIX and DOS machines are available over Internet by anonymous ftp ([sparky.fhcrc.org](http://sparky.fhcrc.org)).

**Constructing blocks databases.** For this work, we began with versions of the Blocks database constructed by PROTOMAT (10) from 504 non-redundant groups of proteins catalogued in PROSITE 8.0 (14) keyed to SWISS-PROT 20 (15). PROTOMAT employs an amino acid substitution matrix at two distinct phases of block construction (16): The MOTIF program uses a substitution matrix when individual sequences are aligned or re-aligned against sequence segments containing a candidate motif (16); The MOTOMAT program uses a substitution matrix when a block is extended to either side of the motif region and when scoring candidate blocks (10). A unitary substitution matrix (matches = 1, mismatches = 0) was used initially, generating 2205 blocks. Next, the BLOSUM program was applied to this database of blocks, clustering at 60%, and the resulting matrix was used with PROTOMAT to construct a second database consisting of 1961 blocks. The BLOSUM program was then applied to this second database, clustering at 60%. This matrix was used to construct version 5.0 of the BLOCKS database from 559 groups in PROSITE 9.00 keyed to SWISS-PROT 22. The BLOSUM program was applied to this final database of 2106 blocks, using a series of clustering percentages to obtain a family of log-odds substitution matrices. This series of matrices is very similar to the series derived from the second database. Approximately similar matrices were also obtained from databases generated by PROTOMAT using the PAM 120 matrix, using a matrix with a clustering percentage of 80% and using just the odd- or even-numbered groups (data not shown).

**Alignments and homology searches.** Global multiple alignments were done using version 3.0 of MULTALIN for DOS computers (17). To provide a positive matrix, each entry was increased by 8 (with default gap penalty of 8). Version 1.6b2 of Pearson's RDF2 program (18) was used to evaluate local pairwise alignments.

Homology searches were done on a Sun Sparcstation using the BLASTP version of BLAST dated 3/18/91 (11) and version 1.6b2 of FASTA (with `ktup=1` and `-o` options) and SSEARCH, an implementation of the Smith-Waterman algorithm (18, 19, 20). The SWISS-PROT 20 databank (15) containing 22,654 protein sequences was searched, and one search was done with each matrix for each of the 504 groups of proteins from PROSITE 8.0. The first of the longest and most distant sequences in the group was chosen as a searching query, inferring

distance from PROTOMAT results and SWISS-PROT names.

In the BLOSUM matrices, the scores for B and Z were made identical to those for D and E, respectively. We used the same gap penalties for all matrices, -12 for the first residue in a gap and -4 for subsequent residues in a gap and -1 for the character X.

The results of each search were analyzed by considering the sequences used by PROTOMAT to construct blocks for the protein group as the true positive sequences and all others as true negatives. BLAST reports the databank matches up to a certain level of statistical significance. Therefore, we counted the number of "misses" as the number of true positive sequences not reported. For FASTA and SSEARCH, we followed the empirical evaluation criteria recommended by Pearson (19); the number of "misses" is the number of true positive scores which ranked below the 99.5th percentile of the true negative scores.

## RESULTS

**Comparison to the Dayhoff matrices.** The BLOSUM series derived from alignments in blocks is fundamentally different from the Dayhoff PAM series which derives from the estimation of mutation rates. Nevertheless, the BLOSUM series based on percent clustering of aligned segments in blocks can be compared to the Dayhoff matrices based on percent accepted mutation (PAM) using a measure of average information per residue pair in bit units called "relative entropy" (9). Relative entropy is zero when the target (or observed) distribution of pair frequencies is the same as the background (or expected) distribution and increases as these two distributions become more distinguishable. Relative entropy was used by Altschul to characterize the Dayhoff matrices, which show a decrease with increasing PAM (9). For the BLOSUM series, relative entropy increases nearly linearly with increasing clustering percentage (Fig. 1). Based on relative entropy, the PAM 250 matrix is comparable to BLOSUM 45 with relative entropy of about 0.4 bit, while PAM 120 is comparable to BLOSUM 80 with relative entropy of about 1 bit. BLOSUM 62 (Fig. 2, lower) is intermediate in both clustering percentage and relative entropy (0.7 bit), and is comparable to PAM 160. Matrices with comparable relative entropies also have similar expected scores.

Some consistent differences are seen when PAM 160 is subtracted from BLOSUM 62 for every matrix entry (Fig. 2, upper). Compared to PAM 160, BLOSUM 62 is less tolerant to substitutions involving hydrophilic amino acids, while it is more tolerant to substitutions involving hydrophobic amino acids. For rare amino acids, especially cysteine and tryptophan, BLOSUM 62 is typically more tolerant to mismatches than is PAM 160.

**Performance in multiple alignment of known structures.** One test of sequence alignment accuracy is to compare the results obtained to alignments seen in 3-D structures. Lipman *et al.* applied a simultaneous multiple alignment program, MSA, to 3 similarly diverged serine proteases of known 3-D structures (21). They found that for 161 closely aligned residue positions, 12 residues were involved in misalignments. We asked how well a hierarchical multiple alignment program, MULTALIN (17), performs on the same proteins using different substitution matrices. Table 1 shows that MULTALIN performs much worse than

MSA using the PAM 120, 160 or 250 matrices, misaligning residues at 30-31 positions. In comparison, MULTALIN with a simple "+6/-1" matrix (that assigns +6 to matches and -1 to mismatches) misaligns residues at 34 positions. In the same test using BLOSUM 45, 62 and 80, MULTALIN misaligned residues at only 6-9 positions. Comparable numbers were obtained when residues that show differences in the positions of side chains were excluded. Therefore BLOSUM matrices produced accurate global alignments of these sequences.

**Performance in searching for homology in sequence databanks.** To determine how BLOSUM matrices perform in databank searches, we first tested them on the G-protein coupled receptors, a particularly challenging group that has been used previously to test searching and alignment programs (10, 18, 22, 23). Three diverse queries, LSHR\$RAT, RTA\$RAT and UL33\$HCMVA, were chosen from among the 114 full-length family members catalogued in PROSITE based on the observation that none detected either of the others in searches. The number of misses was averaged in order to assess the overall searching performance of different matrices for this group. Three different programs were used, BLAST (11) and FASTA (19) and Smith-Waterman (20). BLAST rapidly determines the best ungapped alignments in a databank. FASTA is a heuristic and Smith-Waterman a rigorous local alignment program; both can optimize an alignment by the introduction of gaps. Several BLOSUM and PAM matrices in the entropy range of 0.15 - 1.2 were tested.

Results with each of the 3 programs show that all BLOSUM matrices in the 0.3 to 0.8 range performed better than the best PAM matrix, PAM 200 (Fig. 3). In this range, each BLOSUM matrix missed about 12-25 fewer members than the PAM matrix with similar relative entropy. Therefore, BLOSUM improved detection of members of this family regardless of the searching program used.

To determine whether the superiority of BLOSUM matrices over PAM matrices generalizes to other families, we carried out similar comparative tests for 504 groups of proteins catalogued in PROSITE 8.0. For BLAST, BLOSUM 62 performed slightly better overall than BLOSUM 60 or 70, moderately better than BLOSUM 45 and much better than the best PAM matrix in this test, PAM 140 (Fig. 4). Specifically, BLOSUM 62 was better than PAM 140 for 90 groups, whereas it was worse in only 23 other groups. As a baseline for comparison, we used the simple +6/-1 matrix, which makes no distinction among matches or among mismatches. Compared to +6/-1, BLOSUM 62 performance was better in 157 groups and was worse in 6 groups. Of the 504 groups tested, only 217 showed differences in any comparison. Similar results were obtained for FASTA (data not shown).

Very recently, two updates of the Dayhoff matrices have appeared (24, 25). Both use automated procedures to cluster similar sequences present within an entire protein database, and therefore provide considerably more aligned pairs than were used by Dayhoff. However, in tests of these matrices using BLAST on each of the 504 groups, performance was not noticeably different from that of the Dayhoff PAM 250 matrix which these matrices were intended to replace, much worse than matrices in the BLOSUM series (Fig. 4). Compared to BLOSUM 45, which has similar relative entropy to PAM 250, the matrix of Gonnet, Cohen and Benner (GCB) was worse in 130 groups and better in only 3 groups and the matrix

of Jones, Taylor and Thornton (JTT) was worse in 138 groups and better in only 5 groups.

### **Confirmation of a suspected relationship between transposon ORFs.**

While the above tests demonstrate that BLOSUM matrices perform better overall than PAM matrices, an example indicates the extent to which this improvement can matter in a real situation. We investigated a suspected relationship that is biologically attractive, but is somewhat equivocal when examined by objective criteria. Two groups have noticed a stretch of similarity between the predicted protein from the *Drosophila mauritiana* mariner transposon and that from *Caenorhabditis elegans* transposon Tc1 (S. Emmons, J. Heierhorst, personal communications) (Fig. 5). However, this alignment did not score highly enough to allow its detection in searches using various PAM matrices. In contrast, a BLAST search with BLOSUM 62 using the mariner predicted protein as query detected this alignment as the best in the database (data not shown). An analysis shows non-zero scores taken from the difference matrix of Fig. 1b assigned to each amino acid pair. The higher absolute score for BLOSUM 62 compared to PAM 160 ( $\Sigma=35$  for BL62 > P160 versus  $\Sigma=14$  for BL62 < P160) results from many small differences. When the scores for this alignment were compared to the scores for alignments between one of the sequences and 1000 shuffles of the other, the score using BLOSUM 60 was 7.6 standard deviations (SD) above the mean. In contrast, the score using PAM 160 was only 3.0 SD above the mean with similar results for PAM 250 and PAM 120, accounting for the failure to detect this relationship in previous database searches.

## **DISCUSSION**

We have found that substitution matrices based on amino acid pairs in blocks of aligned protein segments perform better in alignments and homology searches than those based on accepted mutations in closely related groups. Performance was improved overall in every test that we have done, including multiple alignment (MULTALIN), detection of ungapped alignments (BLAST), detection of gapped alignments (FASTA and Smith-Waterman) and determination of the significance of an alignment (RDF2). The importance of such improved performance can be profound for weakly scoring alignments that are not detected in a search or are not trusted. For example, the alignment between predicted proteins encoded by mariner and Tc1 transposons improved by more than 4.5 standard deviations above the mean of comparisons to shuffled sequences when BLOSUM 62 was used instead of PAM matrices.

There are fundamental differences between our approach and that of Dayhoff which could account for the superior performance of BLOSUM matrices in searches and alignments. Dayhoff estimated mutation rates from substitutions observed in closely related proteins and extrapolated those rates to model distant relationships. In our case, frequencies were obtained directly from relationships represented in the blocks, regardless of evolutionary distance. Since blocks were derived primarily from the most highly conserved regions of proteins, it is possible that many of the differences between BLOSUM and PAM matrices arise from different constraints on conserved regions in general. For example, Dayhoff

found asparagine to be the most mutable residue, whereas in blocks, asparagine is involved in substitutions at an average frequency. This could mean that an asparagine located in a mutable region of a protein is itself highly mutable, whereas when it is located in a conserved region, it shows only an average tendency to be involved in substitutions.

Another difference is the larger and more representative dataset used in this work. The Dayhoff frequency table included 36 pairs in which no accepted point mutations occurred. In contrast, the pairs we counted included no fewer than 2369 occurrences of any particular substitution. Scoring differences were especially apparent for pairs involving rare amino acids such as tryptophan and cysteine. Similar findings were made in the two recent updates of the Dayhoff matrix (24, 25). However, in these studies, no evidence was presented that increased data improved performance. Our tests show that the updated Dayhoff matrices still perform poorly overall when compared to BLOSUM 62. This suggests that matrices from aligned segments in blocks, which represent the most highly conserved regions in proteins, are more appropriate for searches and alignments than are matrices derived by extrapolation from mutation rates.

The BLOSUM series depends only upon the identity and composition of groups in PROSITE and the accuracy of the automated PROTOMAT system. While the system itself uses a substitution matrix, iterative application soon leads to nearly the same set of scores, even starting with a unitary matrix or using a representative subset of the groups. Therefore, we do not expect that these substitution matrices will change significantly in the future.

#### **ACKNOWLEDGMENTS**

The suggestion to make a substitution matrix from a blocks database was made by Temple Smith at the 1991 Aspen Center for Physics workshop. We thank Scott Emmons and Jörg Heierhorst for independently pointing out the similarity between mariner and Tc1 predicted proteins, Bill Pearson for advice, and Domokos Vermes for discussions about information theory. This work was supported by a grant from the National Institutes of Health.

## REFERENCES

1. Dayhoff, M. O. & Eck, R. V. (eds.) (1968) *Atlas of protein sequence and structure* (National Biomedical Research Foundation, Silver Spring, Maryland), Vol. 3. (p. 33)
2. McLachlan, A. D. (1971) *J. Mol. Biol.* **61**, 409-424.
3. Feng, D. F., Johnson, M. S. & Doolittle, R. F. (1985) *J. Mol. Evol.* **21**, 112-125.
4. Rao, J. K. M. (1987) *Int. J. Pept. Protein Res.* **29**, 276-281.
5. Risler, J. L., Delorme, M. O., Delacroix, H. & Henaut, A. (1988) *J. Mol. Biol.* **204**, 1019-1029.
6. Smith, R. F. & Smith, T. F. (1990) *Proc. Natl. Acad. Sci. USA* **87**, 118-122.
7. George, D. G., Barker, W. C. & Hunt, L. T. (1990) *Meth. Enzymol.* **183**, 333-351.
8. Dayhoff, M. (1978) *Atlas of protein sequence and structure* (National Biomedical Research Foundation, Washington, D. C.), Vol. 5, suppl. 3. (pp. 345-358)
9. Altschul, S. F. (1991) *J. Mol. Biol.* **219**, 555-565.
10. Henikoff, S. & Henikoff, J. G. (1991) *Nucleic Acids Res.* **19**, 6565-6572.
11. Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. J. (1990) *J. Mol. Biol.* **215**, 403-410.
12. Blahut, R. E. (1987) *Principles and Practice of Information Theory* (Addison-Wesley, Reading, Mass.).
13. Henikoff, S., Wallace, J. C. & Brown, J. P. (1990) *Meth. Enzymol.* **183**, 111-132.
14. Bairoch, A. (1991) *Nucleic Acids Res.* **19**, 2241-2245.
15. Bairoch, A. & Boeckmann, C. (1991) *Nucleic Acids Res.* **19**, 2247-2249.
16. Smith, H. O., Annau, T. M. & Chandrasegaran, S. (1990) *Proc. Natl. Acad. Sci. USA* **87**, 826-830.
17. Corpet, F. (1988) *Nucl. Acids Res.* **16**, 10881-10890.
18. Pearson, W. R. (1990) *Meth. Enzymol.* **183**, 63-98.
19. Pearson, W. R. (1991) *Genomics* **11**, 635-650.
20. Smith, T. F. & Waterman, M. S. (1981) *J. Mol. Biol.* **147**, 195-197.
21. Lipman, D. J., Altschul, S. F. & Kececioglu, J. D. (1989) *Proc. Natl. Acad. Sci. USA* **86**, 4412-4415.
22. Doolittle, R. F. (1990) *Meth. Enzymol.* **183**, 99-110.
23. Attwood, T. K., Eliopoulos, E. E. & Findlay, J. B. C. (1991) *Gene* **98**, 153-159.
24. Gonnet, G. H., Cohen, M. A. & Benner, S. A. (1992) *Science* **256**, 1443-1445.
25. Jones, D. T., Taylor, W. R. & Thornton, J. M. (1992) *CABIOS* **8**, 275-282.
26. Greer, J. (1981) *J. Mol. Biol.* **153**, 1027-1042.

Table 1: Performance of substitution matrices in aligning 3 serine proteases

Matrix	Program	Residue positions missed <sup>1</sup>	
		All positions	Side chains aligned
	MSA	12	6
PAM 120	MULTALIN	31	22
PAM 160	MULTALIN	30	22
PAM 250	MULTALIN	30	22
+6/-1	MULTALIN	34	26
BLOSUM 45	MULTALIN	9	5
BLOSUM 62	MULTALIN	6	4
BLOSUM 80	MULTALIN	9	6

<sup>1</sup>From data of Greer (26), where residues were considered to be aligned whenever  $\alpha$ -carbons occupied comparable positions in space (first column). For a subset (second column), residues were excluded where there were differences in the positions of side chains.

FIG 1. Relationship between percentage clustering and total amino acid pair counts plotted on a log scale (left) and relative entropy (right).

FIG 2. BLOSUM 62 substitution matrix (lower) and the difference matrix (upper) obtained by subtracting the PAM 160 matrix position by position. These matrices have identical relative entropies (0.70); the expected value of BLOSUM 62 is -0.52 and that for PAM 160 is -0.57.

FIG 3. Searching performance of programs using members of the G-protein coupled receptor family as queries and matrices from the BLOSUM and PAM series scaled in half-bits (11). Removal of this family from the BLOCKS database led to a nearly identical matrix with similar performance. Matrices represented (left to right) are BLOSUM 30, 35, 40, 45, 50, 55, 60, 62, 65, 70, 75, 80, 85, 90 and PAM 400, 310, 250, 220, 200, 160, 150, 140, 120, 110, 100. The average number of true positive SWISS-PROT entries missed are shown for LSHR\$RAT, RTA\$RAT and UL33\$HCMVA versus SWISS-PROT 20. Results using BLAST and FASTA or SSEARCH are not comparable to each other, since different detection criteria were used for the three programs.

FIG 4. Searching performance of BLAST using different matrices from the BLOSUM (BL) series, the PAM (P) series and two recent updates of the standard Dayhoff matrix: GCB (24) and JTT (25). Results are based on searches using queries for each of 504 different groups. For each pair of numbers below a box representing a matrix, the first is the number of groups for which BLOSUM 62 missed fewer sequences than that matrix, and the second is the number of groups for which BLOSUM 62 missed more. The vertical distance between each matrix and BLOSUM 62 is proportional to the difference.

FIG 5. Alignment of *D. mauritiana* Mariner predicted protein (aa 245-295) with *C. elegans* TcA (aa 235-285) encoded by Tc1. Difference scores taken from Figure 1b are indicated just below each alignment position. Using RDF2 with BLOSUM 62 for 1000 shuffles and a window size of 10, this alignment scores 64, compared to a mean of 31.4, SD=4.32 for  $z=7.6$ . With PAM 160, the score is 43, compared to a mean of 30.1, SD=4.63 and  $z=3.0$ . With PAM 250,  $z=2.14$  and with PAM 120,  $z=2.98$ .

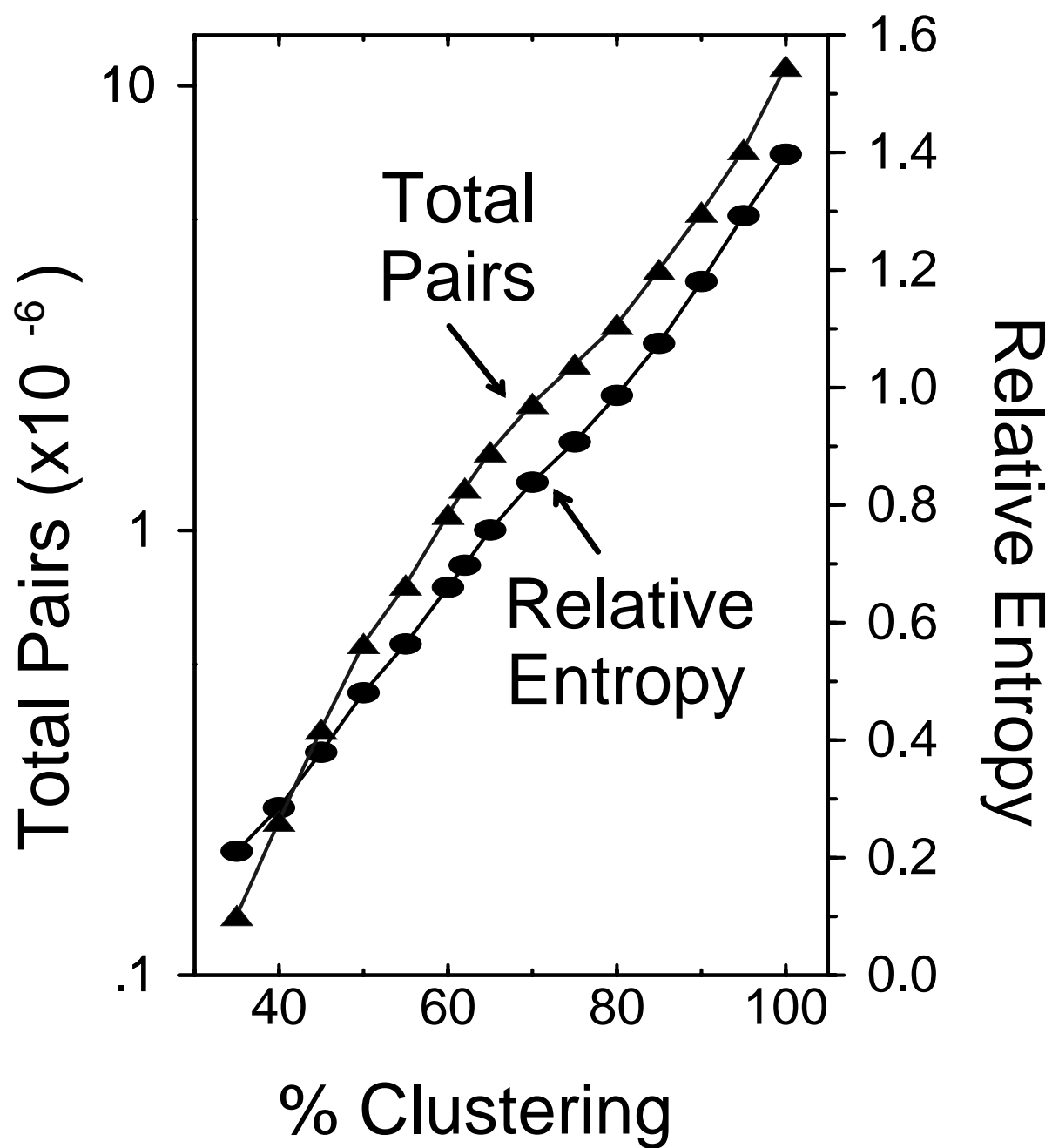


Figure 1

	C	S	T	P	A	G	N	D	E	Q	H	R	K	M	I	L	V	F	Y	W	
	0	-1	1	0	2	1	1	2	1	2	0	0	2	4	1	5	1	2	-2	5	C
		2	0	-2	0	-1	0	0	0	1	0	0	0	1	0	1	-1	1	1	-1	S
C	9		2	-1	-1	-1	0	0	0	0	0	0	-1	0	-1	1	0	1	1	3	T
S	-1	4		2	-2	-1	-1	0	0	-1	-1	-1	1	1	0	-1	0	0	2	1	P
T	-1	1	5		2	-1	-2	-2	-1	0	0	1	1	0	0	1	0	1	1	2	A
P	-3	-1	-1	7		2	0	-1	-2	0	1	1	0	0	-1	0	-1	1	2	4	G
A	0	1	0	-1	4		3	-1	-1	0	0	1	-1	0	-1	0	-1	0	0	0	N
G	-3	0	-2	-2	0	6		2	-1	-1	-1	0	-1	0	0	0	0	2	1	3	D
N	-3	1	0	-2	-2	0	6		1	0	0	2	2	1	-1	0	0	2	2	4	E
D	-3	0	-1	-1	-2	-1	1	6		0	-2	0	1	1	-1	0	0	1	3	3	Q
E	-4	0	-1	-1	-1	-2	0	2	5		2	-1	0	1	0	-1	0	1	2	2	H
Q	-3	0	-1	-1	-1	-2	0	0	2	5		-1	-1	0	-1	1	0	1	3	-4	R
H	-3	-1	-2	-2	-2	-2	1	-1	0	0	8		1	-2	-1	1	1	2	3	1	K
R	-3	-1	-1	-2	-1	-2	0	-2	0	1	0	5		-2	-1	-1	0	1	2	4	M
K	-3	0	-1	-1	-1	-2	0	-1	1	1	-1	2	5		-1	1	0	0	1	3	I
M	-1	-1	-1	-2	-1	-3	-2	-3	-2	0	-2	-1	-1	5		-1	0	-1	1	2	L
I	-1	-2	-1	-3	-1	-4	-3	-3	-3	-3	-3	-3	-3	1	4		0	1	2	4	V
L	-1	-2	-1	-3	-1	-4	-3	-4	-3	-2	-3	-2	-2	2	2	4		-1	-2	1	F
V	-1	-2	0	-2	0	-3	-3	-3	-2	-2	-3	-3	-2	1	3	1	4		-1	2	Y
F	-2	-2	-2	-4	-2	-3	-3	-3	-3	-3	-1	-3	-3	0	0	0	-1	6		-1	W
Y	-2	-2	-2	-3	-2	-3	-2	-3	-2	-1	2	-2	-2	-1	-1	-1	-1	3	7		
W	-2	-3	-2	-4	-3	-2	-4	-4	-3	-2	-2	-3	-3	-1	-3	-2	-3	1	2	11	
	C	S	T	P	A	G	N	D	E	Q	H	R	K	M	I	L	V	F	Y	W	

Figure 2

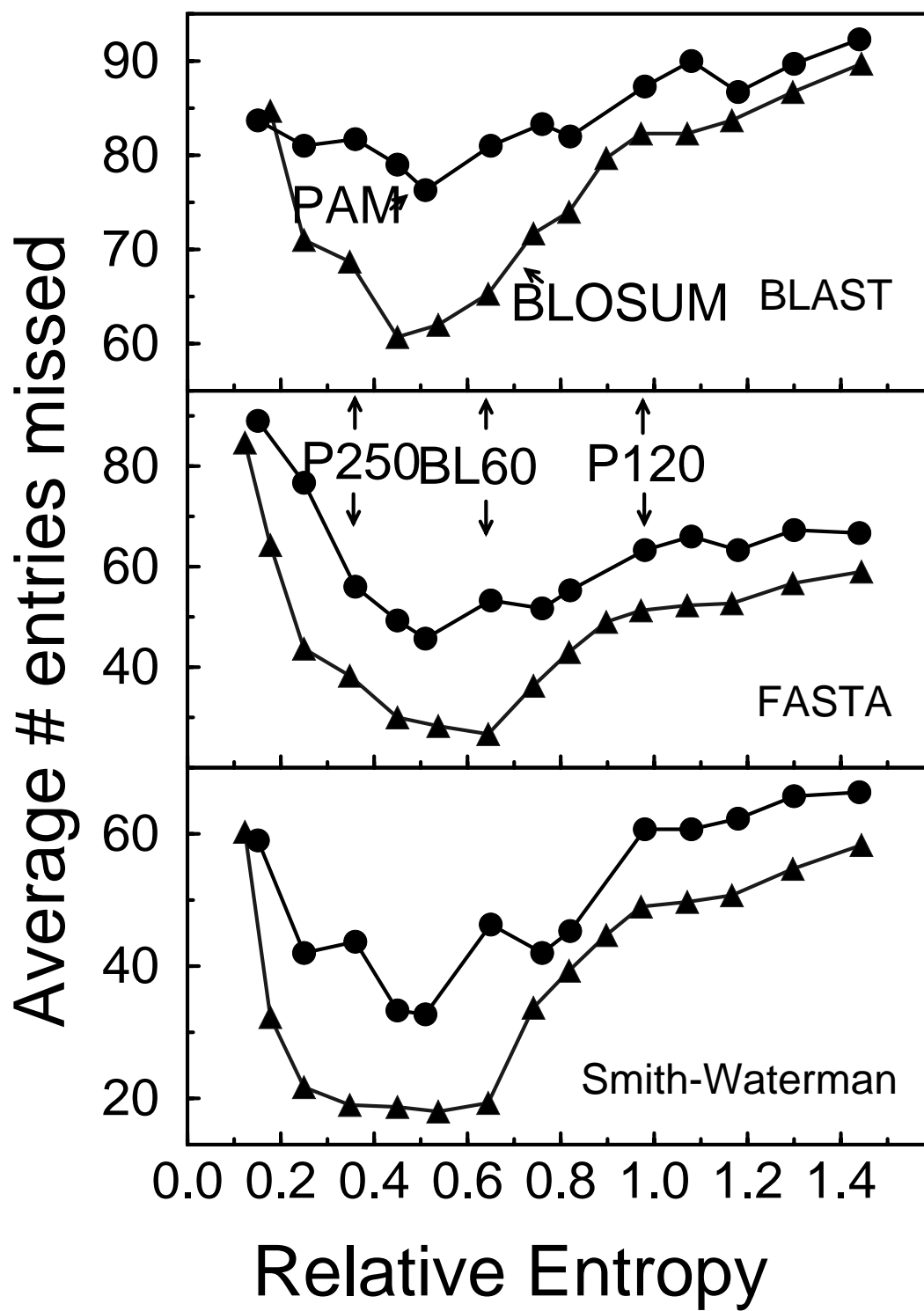


Figure 3

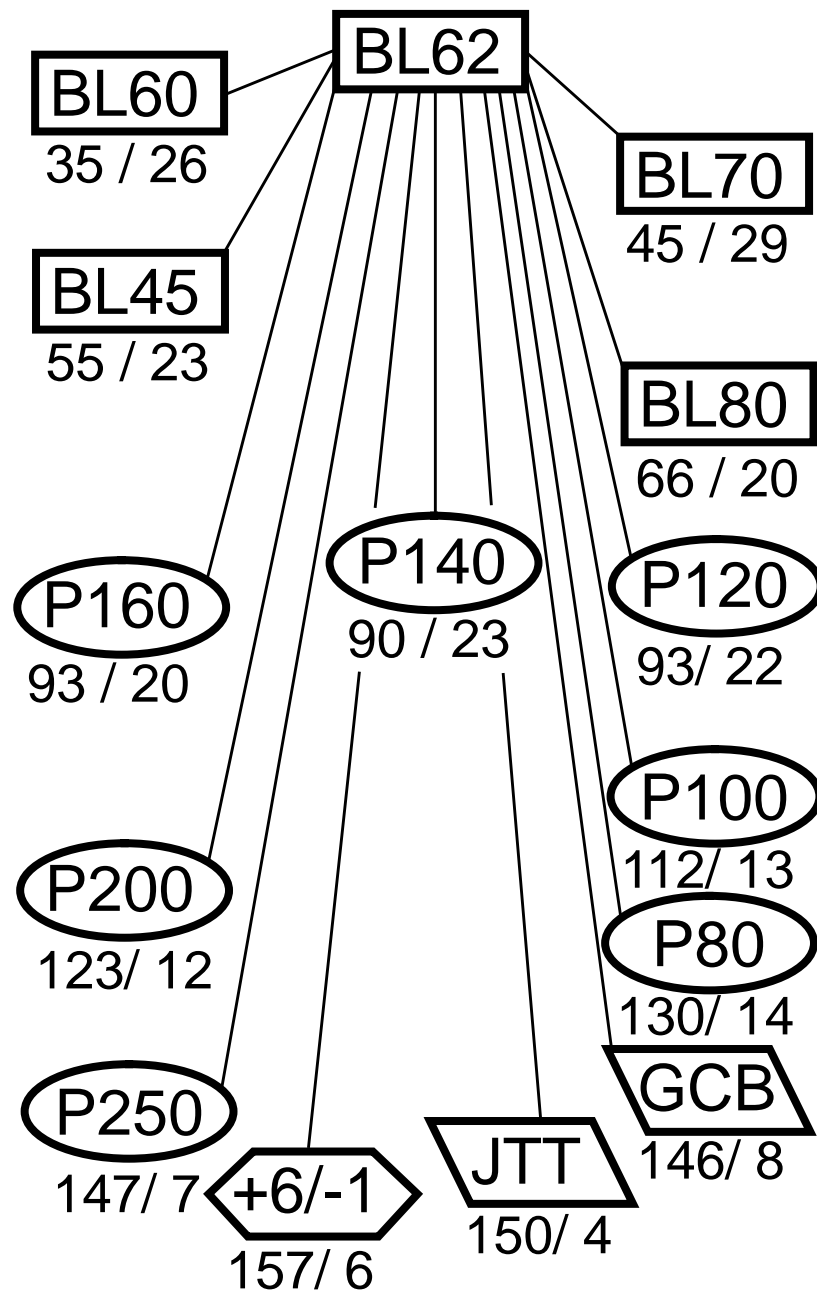


Figure 4

```

Mariner
IFLHDNAPSHTARAVRDTLETNLNWEVLPHAAAYSPDLAPSDY
      :  :: : ::  ::      :      :::: :
Tc1
VFQQDNDPKHTSLHVRSWFQRRHVHLLDWPSQSPDLNPIEH
BL62>P160      23 2 22 1      3  1 4      2  3222  2  2
BL62<P160  1 2  2      1  1      1  2      12  1

```

Figure 5